

Czech Syntactic Analysis Constraint-Based

XDG: One Possible Start

Ondřej Bojar
obo@cuni.cz

Abstract

This article describes an attempt to implement a constraint-based dependency grammar for Czech, a language with rich morphology and free word order, in the formalism Extensible Dependency Grammar (XDG). The grammar rules are automatically inferred from the Prague Dependency Treebank (PDT) and constrain dependency relations, modification frames and word order, including non-projectivity. Although these simple constraints are adequate from the linguistic point of view, their combination is still too weak and allows an exponential number of solutions for a sentence of n words.

1 Motivation

Current approaches to syntactic analysis of Czech and other languages with a high degree of free word order have limitations that are important from the theoretical point of view. First, all the available parsers are restricted to the surface syntactic analysis and there is no simple way of extending them to include a deep syntactic (for instance tectogrammatical) level of representation. Second, the available statistical parsers produce only one solution for a given sentence, ignoring the possibility of the syntactic ambiguity of a sentence. And last but not least, the available parsers¹ are by nature statistical and do not contribute to the explanation of syntactic phenomena very much.

Several declarative (relational) approaches to syntax analysis overcoming these problems are available, including well known formalisms such as HPSG or LFG, or the robust constraint-based dependency parsing by (Foth, Menzel, and Schröder, 2004). Another promising formalism is the Extensible Dependency Grammar² (XDG, (Debusmann et al., 2004)). None of these approaches has ever been tested on a language with rich morphology and freer word order in a large scale.

With respect to ongoing research within the theoretical framework of Functional Generative Description (FGD, see (Sgall, Hajičová, and Panevová, 1986; Hajičová, Sgall, and Partee, 1998)), the Extensible Dependency Grammar is a formalism that excellently fits our needs:

- XDG is dependency based, as FGD is.
- XDG distinguishes between immediate dominance (ID, dependency) relations and linear precedence (LP); constraints are allowed to speak about these two dimensions separately as well as simultaneously and the dimensions are mutually constraining each other. It is easy to handle non-projective constructions in XDG. Both these issues are important with respect to the relatively free word order of Czech.
- XDG allows for handling such dimensions of language description as the deep syntactic (tectogrammatical) level. FGD's main objective is deep syntactic structure.

¹Rare exceptions include an unpublished parser for Czech by Zdeněk Žabokrtský.

²The term *grammar* is used here in the sense of a set of rules underlying syntactic or syntactico-semantic analysis.

- XDG effectively works with ambiguity: morphological, syntactic and lexical ambiguity during parsing is stored in a compact underspecified form as long as possible. The multiplication of orthogonal options is postponed until actually needed.

The task of implementing a large coverage grammar with XDG is interesting for another reason, too. Up to now, only small scale grammars have been implemented in XDG. These grammars illustrated efficient and elegant treatment of various complex syntax and semantic phenomena in XDG (Duchier and Debusmann, 2001; Debusmann and Duchier, 2003). However, the grammars were always tailored to a few test sentences and constraints implemented in XDG never had to cope with syntactic ambiguity of a grammar inferred from a larger amount of data. There are excellent data sources for Czech language from which such a large scale grammar can be collected: the Prague Dependency Treebank (PDT³) and the Czech valency lexicon (Vallex, (Žabokrtský et al., 2002)).

The following text is organized as follows. Section 2.2 summarises the basic principles of constraint parsing and XDG in particular. Section 2.3 offers an overview picture of how different Czech data sources can contribute to an XDG grammar. Section 3 describes the design of my simple Czech grammar for XDG and the procedure of inferring this grammar from the PDT. Results evaluating this grammar on the evaluation part of the PDT are presented in section 4.1 and the most important sources of syntactic ambiguity leading to the enormous number of available solutions are examined in section 4.2.

An outlook of further research and final remarks are given in sections 5 and 6.

2 Introduction

2.1 Properties of Czech Language

Table 1 summarises some of the well known properties of Czech language⁴. Czech is an example of Slavonic languages. It is an inflective language with rich morphology and relatively free word order allowing non-projective constructions. However, there are important word order phenomena restricting the freedom. One of the most prominent examples are clitics, i.e. pronouns and particles that occupy a very specific position within the whole clause. The position of clitics is very rigid and global within the sentence. Locally rigid is the structure of (non-recursive) prepositional phrases or coordination. Other elements, such as the predicate, subject, objects or other modifications may be nearly arbitrarily permuted⁵.

Moreover, like other languages with freer word order, Czech allows non-projective constructions (crossing dependencies). Only about 2% of edges in the PDT are non-projective on the analytic level, but this is enough to make nearly a quarter (23.3%) of all the sentences non-projective.

The task of parsing languages with relatively free word order is much more difficult than parsing of English, for example, and new approaches still have to be searched for. Rich morphology is a factor that makes parsing more time and data demanding.

2.2 Basic principles of XDG

Extensible Dependency Grammar has been designed, as well as a parser implemented by (Duchier and Debusmann, 2001). A grammar in the formalism is specified by providing:

Definition of independent dimensions of linguistic description. Typically, one dimension is used to represent immediate dominance (ID) in the form of (unordered) tree, another dimension is used to represent linear precedence (LP) of the words (ordered projective tree is used on this dimension).

³See (Hajič et al., 2001) and (Hajičová, Panevová, and Sgall, 2000).

⁵Such permutations correspond to the topic-focus articulation of the sentence.

⁵Data by (Collins et al., 1999), (Holan, 2003), Zeman (<http://ckl.mff.cuni.cz/~zeman/projekty/neproji>) and (Bojar, 2003). Consult (Kruijff, 2003) for measuring word order freeness.

| | Czech | English |
|------------------------------|---|--------------------|
| Morphology | rich $\geq 4,000$ tags possible, $\geq 1,400$ seen | limited 50 used |
| Word order | free with rigid global phenomena | rigid |
| Known parsing results | | |
| Edge accuracy | 69.2–82.5% | 91% |
| Sentence correctness | 15.0–30.9% | not reported |

Table 1: Properties of Czech compared to English.

In the Praguian view, only the ID tree is appropriate for “the surface syntactic analysis” of an input sentence, there is no theoretical support for LP trees in FGD. (See Figure 3 below for an illustration.)

In XDG, additional dimensions can be used to represent predicate-argument or semantic structure of the sentence. See (Debusmann and Duchier, 2003) for a more elaborated description of the design of multiple dimensions up to semantic structure. Analogously, the Praguian deep syntactic (tectogrammatical) structure could be added.

General *principles* are used on separate dimensions or pairs of dimensions to enforce properties required on correct analyses. Possible principles include constraints on tree property, projectivity, correspondence between edges on one dimension with edges on another dimension and others.

Lexicon of word forms. The word forms listed in the lexicon are augmented with properties addressing several dimensions of the grammar, the ID and LP dimensions are in general used to express and constrain morphological and syntactic properties of the word forms.

The grammatical properties may be specified not directly with every single word, but rather defined as independent classes. Lexical entries of individual word forms are then described using a boolean combination of the classes. (A class is defined for “word forms being verbal infinitive”, another class is defined for “words accepting direct or indirect object as a modification”. Lexical items of infinitives of ditransitive verbs use the classes: a particular lexical item is “verbal infinitive” and “accepts primary object” and “accepts secondary object”. The same word form can be morphologically ambiguous, so the lexical entry may alternatively belong to the class of pronouns etc.)

The task of syntactic analysis of a sequence of words in XDG is the task of finding tree (graph) structures on all grammar dimensions such that the nodes of the trees correspond one to one to selected lexical entries of the words. The edges of the trees and attributes of the word entries must match all the principles of the grammar. This task is internally expressed as a constraint solving problem (CSP).

The implementation of a XDG parser heavily depends on the constraint solver of the Mozart-Oz programming language⁶. The constraint solver keeps ambiguous information underspecified as long as possible, while using the lexical and grammatical constraints to restrict the set of possible solutions. For example, the remaining possible trees are not directly enumerated, but rather stored implicitly by listing necessary and impossible daughters. If none of the constraints can contribute to the underspecified solution to restrict it, an arbitrary decision (or all possible decisions one by one) is made on one of the still underspecified variables (a lexical entry for a word is chosen, a daughter is said to be present under a node, the value of a morphological property is fixed etc.) Having fixed this particular variable, the constraints obtain a new input and again restrict the available solutions. This is repeated, until a single, fully specified, solution is found, or until a failure (conflicting constraints) is identified.

⁶See <http://www.mozart-oz.org/> for more information and several references.

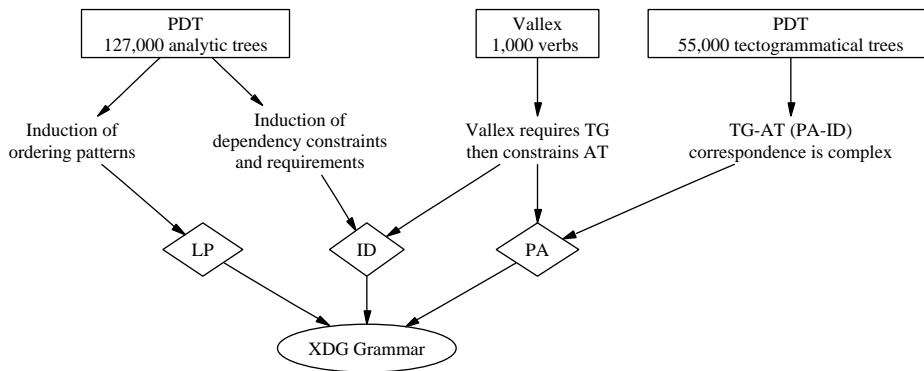


Figure 1: Czech data sources available for XDG grammar.

2.3 Overview of the Intended Multi-dimensional Czech Dependency Grammar

Figure 1 summarises data sources available for a Czech grammar induction. The PDT contains surface syntactic (analytic, AT) as well as deep syntactic (tectogrammatical, TG) sentence annotations. The Czech valency lexicon is under development, and alternatively, the valency lexicon collected while annotating the tectogrammatical level of PDT could be used.

A grammar in the formalism of XDG could be inferred from these sources addressing the immediate dominance (ID), linear precedence (LP) and predicate-argument (PA) dimensions.

Only a part of this overall picture has been implemented so far. First, the correspondence between tectogrammatical and analytic levels is quite complicated, some nodes have to be deleted, some nodes have to be added. Second, the tectogrammatical valency information from Vallex is mostly useful only if a tectogrammatical structure is considered, only then the constraints addressing surface realization can be fully exploited. Therefore, in the first approach the current grammar implementation focuses only on ID and LP levels.

3 Description of the Grammar Parts

The experimental XDG grammar induced from the PDT utilizes basic principles that are linguistically motivated and traditionally used in many varieties of dependency grammars, including XDG. The current XDG grammar extracted from the PDT consists of the following parts: ID Agreement, LP Direction, Simplified ID Frames and ID Look Right. For every part independently, the properties of individual lexical entries (with an arbitrary level of lexicalization) are collected from the training data. The contributions are then combined into XDG lexical entries and classes in a conjunction manner: when parsing, every input word must match one of the observed configurations in all the grammar parts.

For practical reasons (memory and time requirements), the grammar finally used in the XDG parser is restricted to the word forms of the test sentences only. Figure 2 summarizes the pipeline of grammar extraction and evaluation.

3.1 Grammar Non-lexicalized in General

XDG is designed as a lexicalized formalism, most syntactic information is expected to come from the lexicon. Conversely, to make the most use of this approach, the information in an XDG grammar should be as lexicalized as possible.

Despite the size of the PDT (1.5 million tokens), there is not enough data to collect syntactic information for individual word forms and even lemmas.

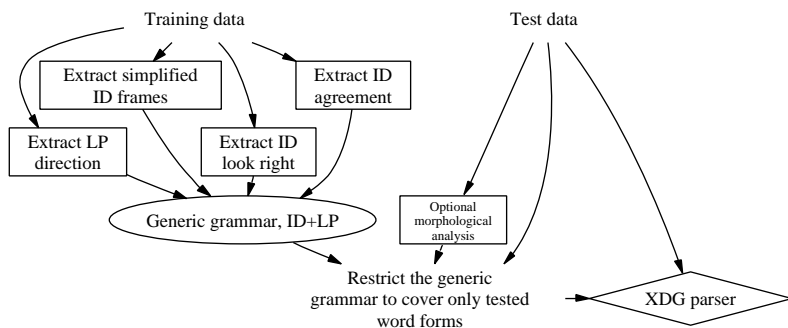


Figure 2: XDG grammar parts and evaluation.

| After having observed | 20,000 | 75,000 | training sentences |
|-----------------------|--------|--------|--------------------|
| lemma (i.e. word) | 1.6 | 1.8 | test sentences |
| full morphological | 110 | 290 | test sentences |
| simplified tag | 280 | 870 | test sentences |

Table 2: Lack of training data in PDT for full lexicalization.

All the grammar parts described below are therefore based on simplified morphological tags only (part and subpart of speech, case, number and gender). Table 2 justifies this simplification. Theoretically, full morphological tags could be used, but we would face sparse data problem if pairs or n -tuples of tags were examined.

3.2 ID Agreement

The ID Agreement part of the grammar allows for a specific edge type between a father and a daughter. The edge type is cross checked in both directions: from the father and from the daughter.

Technically, the lexical entry of a father (with known morphological properties) contains a mapping from edge labels to morphological requirements on any possible daughter. If a daughter is connected via a particular edge label to this father, the daughter’s morphology must match at least one of the requirements. Conversely, the daughter’s lexical entry contains a mapping to restrict the morphology of the father.

This approach helps to reduce morphological ambiguity of nodes: For every node, only such morphological analyses remain allowed which fit the intersection of requirements of all the connected nodes. During parsing, the ambiguous morphology of the node is reduced step by step, as more and more edges are assigned.

3.3 LP Direction

The LP Edge Direction part describes simplified linear precedence rules and handles non-projectivity. In the original design of XDG grammars, motivated by German, the LP dimension is used to describe *topological fields* (Bech, 1955). Unfortunately, the word order of Czech and other Slavonic languages does not exhibit similar word order restrictions in general. (To a very limited extent, one could think about three fields in a clause: preclitic, clitic and postclitic field.) However, there is often an important distinction between dependencies to the left and dependencies to the right.

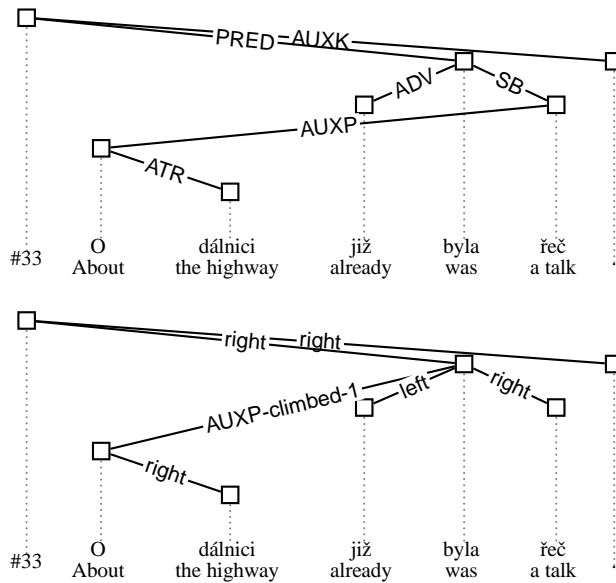


Figure 3: LP dimension to handle edge direction and non-projectivity.

Technically, every father at the LP dimension offers three fields: the left and right fields of unlimited cardinality⁷ and a head field to contain only the father itself. The restriction on specific edge types allowed in a specific direction is handled by another principle: given a daughter connected to the father with an edge of a particular label at the ID dimension, the corresponding LP edge is allowed to have only some of the labels. As illustrated in Figure 3, under the preposition *about*, an (ID) edge labelled ATR can go to the right only, so the corresponding LP edge must have the label RIGHT.

Non-projectivity is forbidden in general but allowed for observed cases.⁸ This constraint is expressed in the LP tree while the ID tree is allowed to be non-projective in general. The LP tree is required to be projective and the exceptions are handled by the so-called *climbing principle*. In order to obtain a projective LP tree from a non-projective one, the tree is “flattened” by climbing. For example, the AUXP edge is non-projective in the ID tree in Figure 3. Moving the corresponding LP edge one step up from the governor *talk* to the governor *was*, the LP edge becomes projective.

To distinguish LP edges that had to climb from LP edges directly corresponding to ID edges, a set of extra LP labels is introduced: AUXP-CLIMBED-1, ATR-CLIMBED-1... The nodes where a climbed edge may land offer not just the left, head and right fields, but also the required amount of specific *-CLIMBED-* edges. There is no restriction on mutual linear ordering of the LEFT/RIGHT and *-CLIMBED-* edges.

The current model still lacks restrictive power to control the clitic position. Similarly, coordination is not modelled properly yet, because the cardinality of left and right fields is unrestricted in general (for example, both members of a coordination are allowed to appear on the same side of the conjunction). More adequate handling of these phenomena remains open for further research.

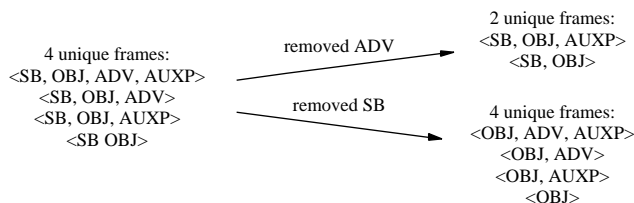
3.4 Simplified ID Frames

One of the crucial principles restricting available sentence analyses in XDG is the valency principle: Every father node allows only specific combinations and cardinalities of outgoing (ID) edges.

⁷In other words, unlimited number of outgoing LP edges can have the label LEFT and all edges labelled LEFT must be present first in the left-to-right ordering of nodes.

⁸Consult (Holan et al., 1998) for a more advanced approach to restricting non-projectivity.

Example: Observed under a verb:



⇒ ADV is more optional than SB.

Figure 4: Identifying optional modifications in order to simplify the set of allowed modification frames.

The Simplified ID Valency Frames ensure that a word doesn't accept implausible combinations of modifiers. Rarely, they ensure that a word has all its "modification requirements" saturated, because most of the modifiers are deletable anyway.

Current approaches⁹ aim at distinguishing *complements* vs. *adjuncts*, i.e. modifications that are typically required vs. optional. However, there is no use of this distinction, if *deletability* of modifications is taken into account (in real Czech sentences, complements are often omitted). Any consistent grammar must reflect this optionality of complements.

The restrictive power of valency frames in XDG should therefore come from *interdependencies* of modifications (e.g. if a secondary object or a specific type of adjunct was observed, a primary object must be present). The set of allowed combinations and cardinalities must be explicitly enumerated in the current XDG implementation. Future versions of this principle might accept a constraint network (for example a set of implications) of interdependencies.

To my knowledge, no published approach aims at discovering such interdependencies of particular modifications so far. On the other hand, there are too many unique frames observed under a given node type, so it is impossible to enumerate all of them.¹⁰

Therefore, I implemented a naive algorithm to infer simplified modification frames: this algorithm automatically simplifies treatment of adjuncts and stores the complexity of interdependencies of other modifications by enumerating them. As sketched in Figure 4, the set of observed modification frames of a specific word class can be simplified by removing different modification types. When an adverbial is removed under a verb, the set of modification frames shrinks to a half in size. When the subject is removed instead, the set does not shrink at all. This indicates that an adverbial has no effect on interdependencies of other modifications: an adverbial may be present or may not—half of the frames was observed with an adverbial, half of the frames had no adverbial.

This simplification is applied iteratively, until the number of unique frames is acceptable. The removed modifications are added to all the frames as optional.

A short example in Figure 5 illustrates the optionality order of modifications observed under (very few) verbs in present tense (POS=V, SUBPOS=B). The most optional modification (AUXP¹¹, a prepositional phrase) is torn off in the first step. The second torn-off modification is an adverbial (ADV), yielding simplified set of modification frames with 36 different frames.

It should be noted that the described solution is by no means a final one. The tasks of inducing modification frames and employing the frames to constrain syntactic analysis are very complex and deserve much deeper research.

⁹See (Sarkar and Zeman, 2000) for comparison and references.

¹⁰Enumerating all seen modification frames would face a severe sparse data problem anyway as the number of unique modification frames steadily grows. In 81,000 sentences, there were 89,000 unique frames observed when describing the frames as lists of simplified tags of all the daughters of a node.

¹¹See (Hajič et al., 2001) for explanation of the labels.

| |
|---|
| Unique observed modification frames: 67 Set sizes when removing specific modifiers: AUXP(50), ADV(50), OBJ(53), SB(55), AUXT(59), AUXG(60), PNOM(61), COORD(62), AUXR(62), AUXY(63), AUXX(65), AUXC(65), APOS(66), HEAD(67), EXD_PA(67), EXD(67) Cumulative simplification: 67→(AUXP)→50→(ADV)→36... |
|---|

Figure 5: Simplifying modifications of verbs in present tense.

3.5 ID Look Right

The generally accepted idea of dependency analysis is that head-daughter dependencies model syntactic analysis best. (Dubey and Keller, 2003) doubt this assumption and document that for German sister-sister dependencies (lexicalized case) are more informative.

| Context used | Neighbours | | Sisters | | |
|--------------|------------|------|---------|------|-------|
| | Head | Left | Right | Left | Right |
| Entropy | 0.65 | 1.20 | 1.08 | 1.14 | 1.15 |

Table 3: Difficulty of predicting an edge label based on simplified tag of a node and a node from close context.

Table 3 gives an indication for Czech: if the structure was already assigned, choosing the edge label is easiest when looking at morphological properties of the node and its head (lowest entropy). Contrary to Dubey and Keller, Czech with a very strong tendency for grammatical agreement confirms the generally accepted view.

The ID Agreement principle is crucial in Czech and it is already employed in the grammar. Table 3 indicates also which context gives the second best hint: the right neighbour, i.e. the following word. Therefore, a new principle was added: ID Look Right: An incoming ID edge to a word must be allowed by the word class of its right neighbour.

The differences among sisters' and neighbours' contributions to the prediction of edge label are not very significant, so adding more constraints of this kind is still under consideration.

4 Results

4.1 Grammar Coverage, Sentence Ambiguity

To evaluate the grammar, only the first fixed point in constraint solving is searched. Given a sentence, the XDG parser propagates all relevant and applicable constraints to reduce the number of analyses and returns an underspecified solution: some nodes may have unambiguously found a governor, for some nodes, several structural assignments may still remain applicable. At the first fixed point, none of the constraints can be used to tell anything more¹².

Two grammars were evaluated: first a version without the Look Right principle, second a version that included the new principle, too. The grammars were trained on sentences from the training part of the PDT and evaluated on 1,800 to 2,000 unseen sentences from the standard evaluation part of the PDT (devtest). The results are displayed in Table 4.

¹²At fixed points, also called choice points, the constraint solver of the underlying system Mozart-Oz makes an arbitrary decision for one of the still underspecified variables and starts propagating constraints again. Other fixed points are reached and eventually a fully specified solution can be printed. Different solutions are obtained by making different decisions at the fixed points. The parser can be instructed to perform a complete search, but in our case there is no point in enumerating so many available solutions.

| Training sentences | 2500 | 5000 | 2500 | 5000 |
|---------------------------|------|------|--------------------------|------|
| Unsolved sentences | | | Avg. ambiguity/node | |
| Without Look Right | 21.1 | 11.9 | 8.09 | 8.91 |
| With Look Right | 25.6 | 15.4 | 8.17 | 9.05 |
| Assigned structural edges | | | Assigned labelled edges | |
| Without Look Right | 4.4 | 3.3 | 3.4 | 2.3 |
| With Look Right | 4.7 | 3.5 | 3.6 | 2.5 |
| Correct structural edges | | | Correctly labelled edges | |
| Without Look Right | 82.3 | 82.5 | 85.9 | 85.9 |
| With Look Right | 81.9 | 81.0 | 85.0 | 83.5 |

Table 4: Results of underspecified solutions.

Note that the number of training sentences was relatively low (around 2 to 5% of the PDT), which explains the relatively high number of unsolved sentences (around 10 to 20%). A wider coverage of the grammar can be easily achieved by training on more data, but this leads to significant growth of the number of solutions available. As indicated in the row Avg. ambiguity/node, a node has 8 to 9 possible governors (regardless the edge label). Compared with the average sentence length of 17.3 words, the grammar reduces the theoretically possible number of structural configurations to a half. At the first fixed point, the parser has enough information to establish only 3 to 5% of edges, an edge with a label can be assigned only to 2 to 4% of nodes. Out of the assigned structural edges, around 82% is correct, out of the assigned labelled edges, around 85% is correct. Again, training on more data should lead to a slightly lower error rate, but significantly less edges securely established, as confirmed by our results.

Contrary to our expectations and preliminary measurements with sentences restricted up to 10 words, the adding the new principle Look Right did not help the analysis. The average ambiguity per node became even higher. There were slightly more edges securely assigned, but the correctness of this assignment has dropped.

4.2 Sources of Ambiguity

Table 5 shows word classes with the highest structural ambiguity at the first fixed point. The data was collected from 1,700 sentences analyzed up to the first fixed point using a grammar trained on 5,000 sentences. The grammar did not use the Look Right principle. The average structural ambiguity is 8.94 possible heads of a node. There were 3.2% edges assigned at the first point, at the level of 82.6% correctness. The average sentence length is 17.32.

The data indicate several problems of the grammar. First, punctuation causes the most severe problem. The reason lies in the fact that the grammar treats punctuation exactly as other word classes and aims at modelling dependency relations of punctuation, although the placement and “dependency assignment” of punctuation depends more on orthography and annotation rules. The fact that the number of possible fathers for an average punctuation mark is higher than the average sentence length can be explained easily: punctuation occurs more often in longer sentences.

Second, coordination is not modelled by the grammar, but it occurs rather often in the corpus. Based on the described constraints, coordinating conjunctions seem highly structurally ambiguous. For example, if the grammar propagated morphological properties of the coordinated members up to the conjunction, the choice of the proper father of the conjunction would become much easier.

Third, word classes without any morphological characteristics such as certain cardinals or adverbs simply “accept many fathers”. The linear precedence constraints would have to be strengthened in order to properly restrict structural ambiguity of these word classes.

| Avg. possible heads | Max. possible heads | Number of observations | Word class |
|---------------------------|---------------------------|---------------------------|-----------------------------------|
| 22.0 | 59 | 2459 | Z : (punctuation) |
| 20.7 | 59 | 1045 | J ^ (conjunctions, coordinating) |
| 19.2 | 49 | 500 | C = (cardinals) |
| 18.9 | 60 | 1035 | D b (adverbs) |
| 18.2 | 57 | 543 | J , (conjunctions, subordinating) |
| 17.3 | 59 | 2611 | R R (prepositions) |
| 17.2 | 44 | 448 | D g (adverbs) |
| 15.6 | 38 | 130 | T T (particles) |
| 13.9 | 47 | 1568 | V B (verbs) |
| 12.8 | 36 | 519 | V f (verbs) |
| 12.6 | 37 | 201 | R V (prepositions) |
| 12.1 | 49 | 8949 | N N (nouns) |
| 11.1 | 39 | 1197 | V p (verbs) |
| 9.6 | 37 | 535 | P 7 (pronouns) |
| 9.1 | 26 | 166 | V s (verbs) |
| | | ... | |
| 6.1 | 49 | 3310 | A A (adjectives) |

Table 5: Sources of structural ambiguity.

Fourth, clause structure has to be modelled. Modelling clause structure would help not just to reduce the structural ambiguity of subordinating conjunctions and verbs, but also to partition the input sequence into separate chunks in which constraints would be applied faster.

The “traditional” reason for syntactic ambiguity, i.e. noun phrase or prepositional phrase attachment, comes at earliest in the fifth place. This problem is still quite severe: for instance 12.1 possible fathers of a noun in a sentence.

Several factors contribute to the structural ambiguity at the first fixed point. One of the very important factors is the lack of probabilities. For example, genitive noun phrases following nouns depend on the preceding noun with a very high probability. However, there are certain verbs that require (or at least accept) genitive complements. Because the grammar is not lexicalized yet, all verbs seem to offer an edge for the genitive noun phrases. With probabilities, the parser could assign the genitive noun phrases to the preceding noun without making a non-frequent error, even if the parser would not use lexicalized information yet.

5 Discussion and Further Research

The presented results indicate several weak points in the described approach to constraint-based dependency parsing. All these points remain open for further research.

First, the tested version of XDG parser could not make any use of frequency information contained in the PDT.¹³ (Dienes, Koller, and Kuhlmann, 2003) attempt at guiding the XDG parser by frequency information, but the research is still in progress. A similar constraint-based dependency parsing by (Heinecke et al., 1998) inherently includes weight of constraints, but no directly comparable results were published so far. ((Foth, Menzel, and Schröder, 2004) report edge accuracy of 96.63% on a corpus of 200 sentences with average length 8.8 words.)

Second, the current grammar relies on very few types of constraints. More constraints of different kinds have to be added to achieve better propagation. A related problem is the locality of the constraints.

¹³In an experiment, frequency information was used as a threshold to ignore rare edge assignments. The thresholding resulted in lower coverage *and* lower precision.

All the current constraints rely on a too local context. There are too many analyses available, because the local constraints are not powerful enough to check invariant properties of clauses or sentences as a whole.

Third, there are several kinds of expressions that in fact have no dependency structure, such as names, dates and other multi-word expressions. The “dependency” analysis of such expressions in the PDT reflects more the annotation guidelines than some linguistic motivation. Separate treatment of these expressions by means of a sub-grammar would definitely improve the overall accuracy.

6 Conclusion

I described an experiment with constraint based dependency parsing of a language with rich morphology and freer word order. Although the constraints are linguistically adequate and serve well when employed on small-scale corpora, they face a serious problem when trained on large data sets. The constraints are too local and weak in order to restrict the number of available solutions.

To amend this problem, new kinds of constraints have to be developed. In order to achieve a plausible solution quickly, some sort of probabilistic guidance must be added to the constraint solver.

7 Acknowledgement

I'm grateful to Ralph Debusmann for his explanatory and immediate implementation support of new features needed in the XDG parsing system for this experiment. The work could not have been performed without the support of Programming Systems Lab headed by Gert Smolka (Universität des Saarlandes) and without the insightful guidance by Geert-Jan Kruijff and Denys Duchier.

References

- Bech, Gunnar. 1955. *Studien über das deutsche Verbum infinitum*. 2nd unrevised edition published 1983 by Max Niemeyer Verlag, Tübingen (Linguistische Arbeiten 139).
- Bojar, Ondřej. 2003. Towards Automatic Extraction of Verb Frames. *Prague Bulletin of Mathematical Linguistics*, (79–80):101–120.
- Collins, Michael, Jan Hajič, Eric Brill, Lance Ramshaw, and Christoph Tillmann. 1999. A Statistical Parser of Czech. In *Proceedings of 37th ACL Conference*, pages 505–512, University of Maryland, College Park, USA.
- Debusmann, Ralph and Denys Duchier. 2003. A meta-grammatical framework for dependency grammar.
- Debusmann, Ralph, Denys Duchier, Alexander Koller, Marco Kuhlmann, Gert Smolka, and Stefan Thater. 2004. A relational syntax-semantics interface based on dependency grammar. Technical report. Available at <http://www.ps.uni-sb.de/Papers>.
- Dienes, Péter, Alexander Koller, and Marco Kuhlmann. 2003. Statistical a-star dependency parsing. In Denys Duchier, editor, *Prospects and Advances of the Syntax/Semantics Interface*, pages 85–89, Nancy.
- Dubey, Amit and Frank Keller. 2003. Probabilistic parsing for German using sister-head dependencies. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 96–103, Sapporo.
- Duchier, Denys and Ralph Debusmann. 2001. Topological dependency trees: A constraint-based account of linear precedence. In *39th Annual Meeting of the Association for Computational Linguistics (ACL 2001)*.
- Foth, Kilian, Wolfgang Menzel, and Ingo Schröder. 2004. Robust parsing with weighted constraints. *Natural Language Engineering*. in press.
- Hajič, Jan, Jarmila Panevová, Eva Buráňová, Zdeňka Urešová, and Alla Bémová. 2001. A Manual for Analytic Layer Tagging of the Prague Dependency Treebank. Technical Report TR-2001-, ÚFAL MFF UK, Prague, Czech Republic. English translation of the original Czech version.
- Hajičová, Eva, Jarmila Panevová, and Petr Sgall. 2000. A Manual for Tectogrammatic Tagging of the Prague Dependency Treebank. Technical Report TR-2000-09, ÚFAL MFF UK, Prague, Czech Republic. In Czech.
- Hajičová, Eva, Petr Sgall, and Barbara Partee. 1998. *Topic-focus articulation, tripartite structures, and semantic content*. Kluwer, Dordrecht, ISBN 0-7923-5289-0.
- Heinecke, Johannes, Jürgen Kunze, Wolfgang Menzel, and Ingo Schröder. 1998. Eliminative parsing with graded constraints. In *Proceedings of COLING-ACL Conference*, Montreal, Canada.
- Holan, T., V. Kuboň, K. Oliva, and M. Plátek. 1998. Two Useful Measures of Word Order Complexity. In A. Polguere and S. Kahane, editors, *Proceedings of the Coling '98 Workshop: Processing of Dependency-Based Grammars*, Montreal. University of Montreal.
- Holan, Tomáš. 2003. K syntaktické analýze českých(!) vět. In *MIS 2003*. MATFYZPRESS, January 18–25, 2003.
- Kruijff, Geert-Jan M. 2003. 3-phase grammar learning. In *Proceedings of the Workshop on Ideas and Strategies for Multilingual Grammar Development*.
- Sarkar, Anoop and Daniel Zeman. 2000. Automatic Extraction of Subcategorization Frames for Czech. In *Proceedings of the 18th International Conference on Computational Linguistics (Coling 2000)*, Saarbrücken, Germany. Universität des Saarlandes.
- Sgall, Petr, Eva Hajičová, and Jarmila Panevová. 1986. *The Meaning of the Sentence and Its Semantic and Pragmatic Aspects*. Academia/Reidel Publishing Company, Prague, Czech Republic/Dordrecht, Netherlands.
- Žabokrtský, Zdeněk, Václava Benešová, Markéta Lopatková, and Karolina Skwarská. 2002. Tektogramaticky anotovaný valenční slovník českých sloves. Technical Report TR-2002-15, ÚFAL/CKL, Prague, Czech Republic.